

Machine Learning v5. Neural Network Learning.

1. Cost Function.

• neural network for classification.

• L : total # of layers.

• s_ℓ : # of units in layer ℓ .

(1) Binary classification.

$y=0$ or 1 . ; 1 output unit. $h_{\Theta}(x) \in \mathbb{R}$.

$s_L=1$, $K=1$.

(2) Multi-class classification.

$y \in \mathbb{R}^K$. e.g. $\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$, $\begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$,

K output units.

$h_{\Theta}(x) \in \mathbb{R}^K$. $s_L=K$. ($K \geq 3$)

Cost function

Logistic regression. $J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\Theta}(x^{(i)}) + (1-y^{(i)}) \cdot \log (1-h_{\Theta}(x^{(i)})) \right]$

↓ generalized.

$$+ \underbrace{\frac{\lambda}{m} \sum_{j=1}^K \theta_j^2}_{\text{regularization.}}$$

Neural network. $h_{\Theta}(x) \in \mathbb{R}^K$; $(h_{\Theta}(x))_{(i)} = i^{\text{th}}$ output.

$$J_{\Theta}(x) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \cdot \log (h_{\Theta}(x^{(i)})_k) + (1-y_k^{(i)}) \log (1-h_{\Theta}(x^{(i)})_k) \right]$$

←
sum over output
units

$$+ \frac{\lambda}{2m} \sum_{i=1}^{L-1} \sum_{\ell=1}^{s_\ell} \sum_{j=1}^{s_{\ell+1}} (\theta_{j\ell}^{(i)})^2$$

2. Back Propagation Algorithm

min $J(\Theta)$ \rightarrow need to compute $J(\Theta)$

$$\left[-\frac{\partial}{\partial \Theta_{ij}} J(\Theta) \right]$$

(1) Gradient computation

Consider only a single training example (x, y)

Forward propagation.

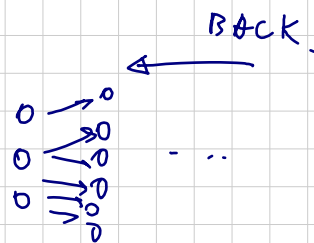
$$- a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$- a^{(2)} = g(z^{(2)}) \quad (\text{add } a^{(2)})$$

\vdots

$$- a^{(4)} = h_{\Theta}(x) = g(z^{(4)})$$



Back propagation

intuition. $f_j^{(l)}$ = "error" of node j in layer l .

recall: $a_j^{(l)}$ activation unit in \dots

for each output unit. (layer $L=4$)

$$f_j^{(4)} = a_j^{(4)} - y_j \quad \text{"labelled" training set.}$$

vectorized. $\left(\begin{array}{l} f^{(4)} = \underline{a}^{(4)} - \underline{y} \end{array} \right)$ $\in \mathbb{R}^K$ \leftarrow # of output unit.

back. $\left(\begin{array}{l} f^{(3)} = (\Theta^{(3)})^T f^{(4)} \cdot \left[g'(z^{(3)}) \right] \end{array} \right)$

$$f^{(2)} = (\Theta^{(2)})^T \cdot f^{(3)} \cdot \left[g'(z^{(2)}) \right] \quad \underline{a}^{(2)} \cdot (1 - \underline{a}^{(2)})$$

$$g = \frac{1}{1 + e^{-x}}$$

It can be shown that:

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \cdot f_i^{(l+1)} \quad (\text{for } j \neq 0)$$

Backpropagation Algorithm

Training set $\{ (x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)}) \}$

Set $\Delta_{ij}^{(l)} = 0$ for all l, i, j \rightarrow used to compute $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

For $i = 1:m$ % loop thru. training set.

- ① set $a^{(1)} = x^{(i)} \rightarrow$ input for i^{th} training example.
- ② perform forward propagation to compute $a^{(l)}$ for $l = 1, 2, 3 \dots L$
- ③ using $y^{(i)}$, $f^{(L)} = a^{(L)} - y^{(i)}$
 $f^{(L-1)}, f^{(L-2)} \dots$

$$\Delta_{ij}^L = \Delta_{ij}^L + a_j^{(L)} \cdot f_i^{(L+1)}$$

\downarrow vect.

$$\Delta^L = \Delta^L + f^{(L+1)} (a^{(L)})^T$$

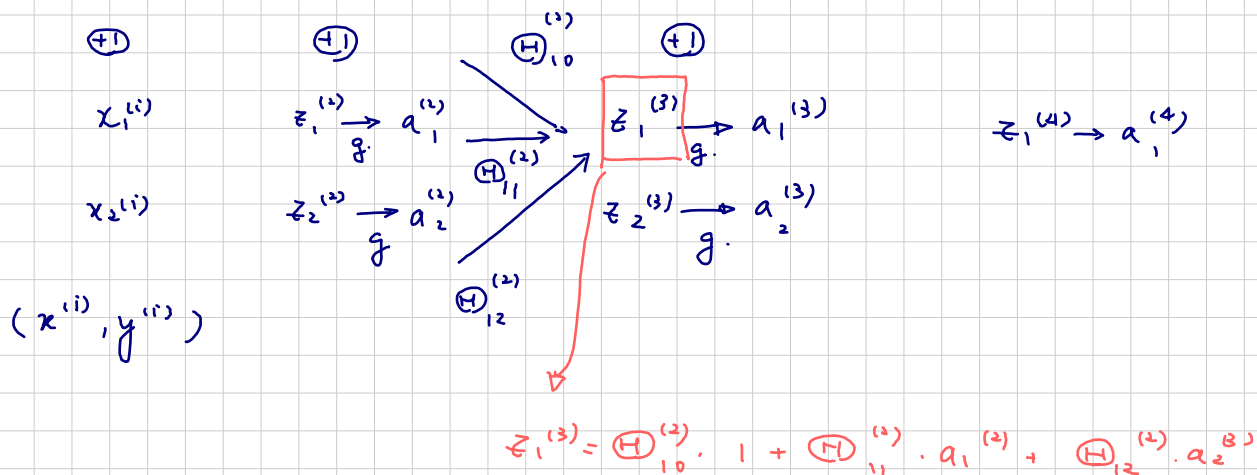
$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \quad \text{if } j \neq 0$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j = 0$$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$$

3. Backpropagation Intuition.

(1) Forward propagation.



Now, notice that backpropagation is doing something very similar to this, but in the reversed dir. (\leftarrow) instead of (\rightarrow)

(2) What is back propagation doing?

$$J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \cdot \log[h_{\Theta}(x^{(i)})] + (1 - y^{(i)}) \cdot \log[1 - h_{\Theta}(x^{(i)})] \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} [\theta_{ji}^{(l)}]^2$$

(for now, assume $\lambda=0$, ignore regularization)

• focus on single example $(x^{(i)}, y^{(i)})$, case of 1 output unit.

• consider cost: $\text{cost}^{(i)} = y^{(i)} \cdot \log h_{\Theta}(x^{(i)}) + (1 - y^{(i)}) \cdot \log (1 - h_{\Theta}(x^{(i)}))$

$$\text{cost}^{(i)} \approx (h_{\Theta}(x^{(i)}) - y^{(i)})^2$$

(how close is $y^{(i)}$ to output $h_{\Theta}(x^{(i)})$?)

Now, $\delta_j^{(l)} =$ "error" of cost for $a_j^{(l)}$ (unit j in layer l)

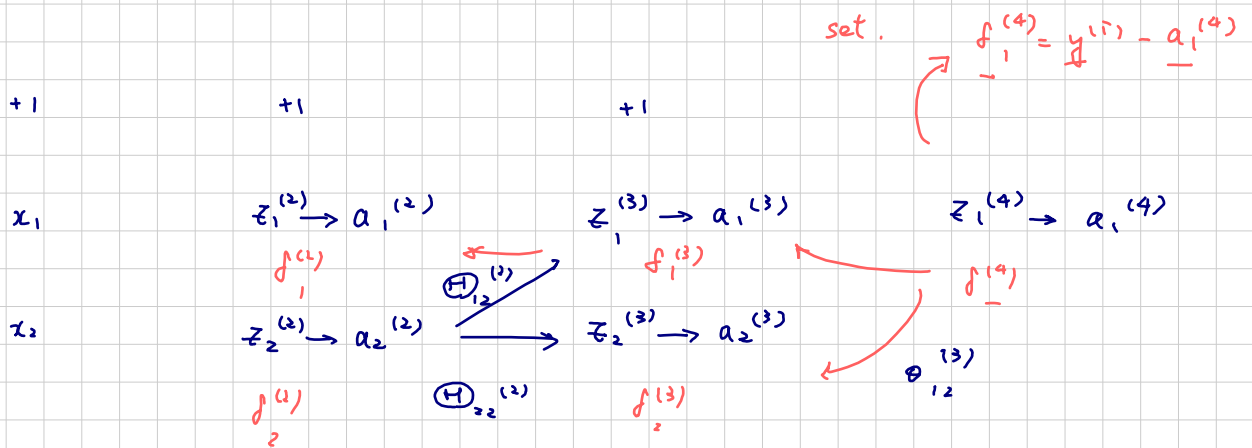
Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}^{(i)}$ for $j > 0$. Where.

$$\text{cost}^{(i)} = y^{(i)} \log h_{\oplus}(x^{(i)}) + (1 - y^{(i)}) \cdot \log h_{\ominus}(x^{(i)})$$

(measure of how much we need to change $z_j^{(l)}$ to affect $h_{\oplus}(x^{(i)})$).

[partial derivatives of cost func. wrt. intermediate values,

so as to affect the output].



Q

How did we arrive at $\delta_2^{(2)}$? "Weighted sum of the activation unit in L1. layer"

$$\delta_2^{(2)} = \Theta_{12}^{(2)} \delta_1^{(3)} + \Theta_{22}^{(2)} \delta_2^{(3)}$$

$$\text{Q. } \delta_2^{(3)} = \Theta_{12}^{(3)} \cdot \delta_1^{(4)}$$

Machine Learning Week 5 Pt. 2: Backpropagation in Practice.

4. Implementation Note: Unrolling Parameters.

Advanced optimization

function. [val, gradient] = costFunction(theta)
... $\xrightarrow{\mathbb{R}^{n+1}}$ $\xrightarrow{\mathbb{R}^{k+1} \text{ (vectors)}}$
optTheta = fminunc(@costFunction, initialTheta, options)

Neural Network ($l=4$) \rightarrow not "vectors" anymore, instead "matrices"
 $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices ($\Theta_1, \Theta_2, \Theta_3$)
 $D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D_1, D_2, D_3)
 \hookrightarrow unroll into vectors.

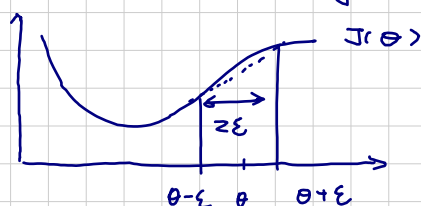
• reshape. \rightarrow slice elements and fit into matrix.

Learning Algorithm.

• Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)} \rightarrow$ unroll to get initialTheta.

5. Gradient Checking

- eliminates error from escaping verification. (BP is error-prone but sometimes not detectable)
- Numerical estimation of gradients.



- get secant line b/w $(\theta-\epsilon, \theta+\epsilon)$

$$\frac{dJ(\theta)}{d\theta} \approx \frac{J(\theta+\epsilon) - J(\theta-\epsilon)}{2\epsilon}, \quad \epsilon \sim 10^{-4}.$$

$$\left(\text{ult. } \frac{dJ(\theta)}{d\theta} = \frac{J(\epsilon+\theta) - J(\theta)}{\epsilon} \right)$$

When, $\theta \in \mathbb{R}^n$. Unrolled version of $\Theta^{(1)}$, $\Theta^{(2)}$, $\Theta^{(3)}$]

so $\theta = (\theta_1, \theta_2, \dots, \theta_n)$

$$\begin{pmatrix} \frac{\partial}{\partial \theta_1} J(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_n} J(\theta) \end{pmatrix} \approx \frac{J(\theta_1 + \epsilon, \theta_2, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \dots, \theta_n)}{2\epsilon}$$

change only 1. and hold other const.

approximating $\nabla J(\theta)$

from backprop.

key: check grad \approx DVec

Note

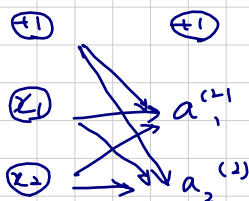
- ① BP to compute DVec. [unrolled $D^{(1)}$, $D^{(2)}$, $D^{(3)}$]
 - ② Implement numerical grad. check. to compute grad approx.
 - ③ Make sure they give similar values.
 - (*) ④ Turn off gradient checking, use backprop. for learning. [D, $f^{(1)}$, $f^{(2)}$...]
- very expensive.
- Remember! Disable gradient checking else code will run slow

6. Random Initialization.

We need initial value for Θ

- set initialTheta = zeros(n, 1) \rightarrow logster \checkmark

- NN.



$$\Theta_{ij}^{(l)} = 0 \quad \forall l, i, j$$

After each update, weight pairs are still same (diff. value) btwn. each other.

$$a_1^{(2)} = a_2^{(2)}, \text{ also } f_1^{(2)} = f_2^{(2)}$$

$$\Rightarrow \frac{\partial}{\partial \Theta_{o_1}^{(1)}} J(\theta) = \frac{\partial}{\partial \Theta_{o_2}^{(1)}} J(\theta) \Rightarrow \Theta_{o_1}^{(1)} = \Theta_{o_2}^{(1)}$$

⇒ Random Initialization: symmetry breaking.

Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$

(i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)

no (same ϵ in
prev. sect.)

7. Summary.

(1) Training a neural network.

- pick a network architecture (connectivity)

- layers, # of act. units.

< input unit: dimension of features $x^{(i)}$.
output .. : # of classes.

multiclass → $y = \{1, 2, 3, \dots, 10\}$.

$$y = \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \text{or} \quad \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$[y=1]$ $[y=2]$

- Reasonable default.

$\left\{ \begin{array}{l} 1 \text{ hidd. layer} \\ > 1 \text{ hidd. layer, w/ same \# of hidd. units.} \end{array} \right.$

(the more the better) ⇔ complex? ✓

1) Randomly initialize weights

2) FP: $h_{\Theta}(x^{(i)})$ $\neq x^{(i)}$

3) $J(\Theta)$

4) BP: $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$

- for loop over examples (get $n^{(l)}$ and $f^{(l)}$)

$l=2, \dots, L$)

5) Gradient checking: comp. $\frac{\partial J(\Theta)}{\partial \theta_{ik}}$ BP

num. est. of $\frac{\partial J(\Theta)}{\partial \theta}$

disable

6) $\left\langle \begin{array}{l} \text{Grad. descent} \\ \text{Adv. opt. method.} \end{array} \right\rangle$ BP to min. $J(\Theta)$ over Θ .

- $J(\Theta)$ non-convex \Rightarrow can get stuck in local minima.

\hookrightarrow good enough in practise.